# Predicting NFL Running Back Yardage Based on Prior Season's Performance



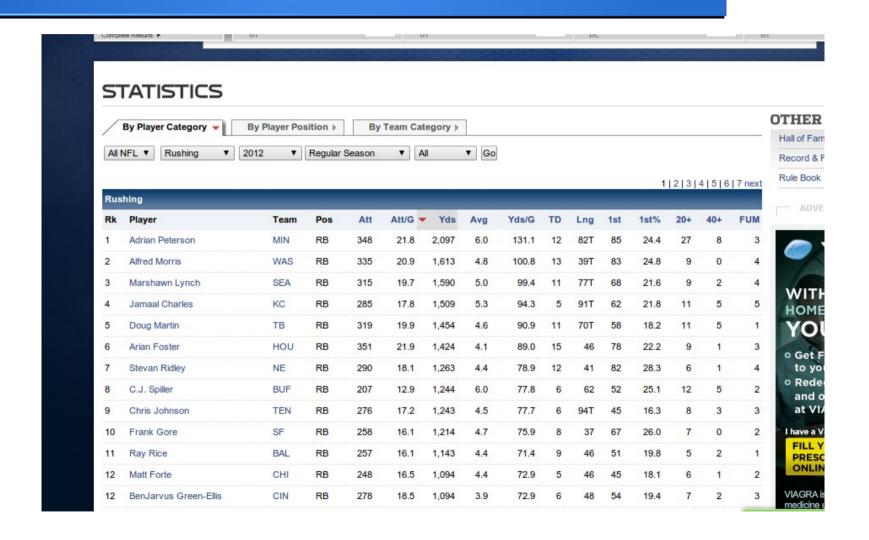Zachary Hogan | General Assembly | Data Science | Spring 2013

# The Problem & Hypothesis

- Predict NFL running back yardage

- Analysis will use past season's data to build a model

- Hypothesis: using running back performance measures as features in a linear model will produce an accurate prediction of future performance
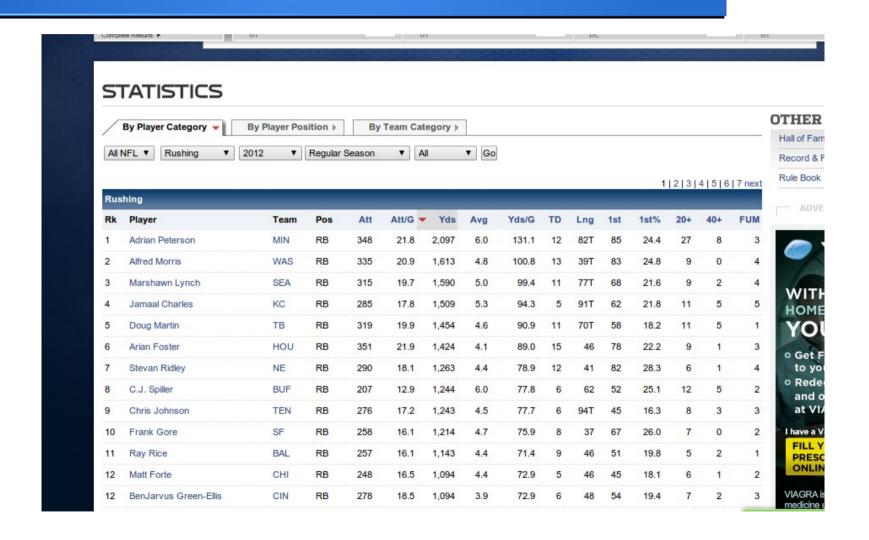
# The Data

- Abundance of potential web-sites to collect NFL data

- Not very many unified data sets

- NFL.com most features in one place, trusted source

- 2011 data - training set

- 2010 data - test set

- Will use the model on the 2012 data

# The Data – Scraping

- Inspect the HTML

- Data is contained within single table element

- Multiple pages for each year of data

- Store each of these URLs in an array

Zachary Hogan | General Assembly | Data Science | Spring 2013

# The Data – Scraping

# The Data – Scraping

- Inspect the HTML

- Data is contained within single table element

- Multiple pages for each year of data

- Store each of these URLs in an array

# The Data – Scraping

# The Data – Scraping

- Inspect the HTML

- Data is contained within single table element

- Multiple pages for each year of data

- Store each of these URLs in an array

# The Data – Cleaning

- Multiple columns contain numerous escape characters ( \n, \t)

- Longest rush column contains 'T' characters

- Commas from the total yardage column

- 2010 test data contains escape characters in additional columns

# The Data – Storage

- Merge future season yardage to the test and training sets

- In Python, store the data in Pandas data frame

- Export as csv, to transfer to R or anywhere else we might need it

- Left with a Test set, a Training set, and the data we will predict on

# Statistical Method

|  | Continuous | Categorical |
|---|---|---|
| Supervised | regression | classification |
| Unsupervised | dimension reduction | clustering |

# Statistical Method

| | Continuous | Categorical |
|---|---|---|
| Supervised | **REGRESSION** | classification |
| Unsupervised | dimension reduction | clustering |

# Statistical Method

- Multivariate Linear Regression

- 2011 data is the training set

- 2010 will be our test set

- If the model performs well on the 2010 data, we will use it to predict the 2013 results based on the 2012 data

# Statistical Method

```
train_fit <- lm(yards_2012 ~ att + att.game +
yards + avg + yards.game + TD + long + X1st
+ X20. + X40. + fumbles, data=train_data11)
```

# Statistical Method

- This fit produces an $R^2$ of 0.6006

- Attempt backwards elimination

# Statistical Method

```
train_data10["predict_yds"] <-
predict(train_fit, train_data10)
```

# Statistical Method 2010 Test Set

| Running Back | 2011 Predicted Yards | 2011 Actual Yards | Difference |
|---|---|---|---|
| Arian Foster | 1173 | 1224 | 51 |
| Jamaal Charles | 564 | 83 | -481 |
| Michael Turner | 1131 | 1340 | 209 |
| Chris Johnson | 1303 | 1047 | -256 |
| Maurice Jones-Drew | 793 | 1606 | 813 |
| Adrian Peterson | 988 | 970 | -18 |
| Rashard Mendenhall | 1367 | 928 | -441 |
| Steven Jackson | 890 | 1145 | 255 |
| Ahmad Bradshaw | 988 | 659 | -329 |
| Ray Rice | 777 | 1364 | 587 |

# Statistical Method
# 2010 Test Set

- Root Mean Squared Error – 231.8 yards

- Use this to measure to compare the accuracy of future modifications to the model

# Statistical Method 2013 Predictions

| Running Back | 2012 Yards | 2013 Predicted Yardage |
|---|---|---|
| Adrian Peterson | 2097 | 1700 |
| Alfred Morris | 1613 | 1152 |
| Marshawn Lynch | 1590 | 1022 |
| Jamaal Charles | 1509 | 619 |
| Doug Martin | 1454 | 1161 |
| Arian Foster | 1424 | 1218 |
| Stevan Ridley | 1263 | 670 |
| C.J. Spiller | 1244 | 732 |
| Chris Johnson | 1243 | 746 |
| Frank Gore | 1214 | 743 |

# Business Applications

- NFL team management

- NFL television analysis

- Fantasy Football

# Conclusion

- Reasonably accurate prediction on the 2010 test data

- Lots of room for improvement

    - Use only data from top running backs, combine it with multiple years of data

    - Additional features

    - Predict rank rather than yardage

# Conclusion

Thank you!